

# Tracking Initiative in Collaborative Dialogue Interactions

Jennifer Chu-Carroll and Michael K. Brown

Bell Laboratories

Lucent Technologies

600 Mountain Avenue

Murray Hill, NJ 07974, U.S.A.

E-mail: {jenc, mkb}@bell-labs.com

## Abstract

In this paper, we argue for the need to distinguish between *task* and *dialogue* initiatives, and present a model for tracking shifts in both types of initiatives in dialogue interactions. Our model predicts the initiative holders in the next dialogue turn based on the current initiative holders and the effect that observed cues have on changing them. Our evaluation across various corpora shows that the use of cues consistently improves the accuracy in the system's prediction of task and dialogue initiative holders by 2-4 and 8-13 percentage points, respectively, thus illustrating the generality of our model.

## 1 Introduction

Naturally-occurring collaborative dialogues are very rarely, if ever, one-sided. Instead, initiative of the interaction shifts among participants in a primarily principled fashion, signaled by features such as linguistic cues, prosodic cues and, in face-to-face interactions, eye gaze and gestures. Thus, for a dialogue system to interact with its user in a natural and coherent manner, it must recognize the user's cues for initiative shifts and provide appropriate cues in its responses to user utterances.

Previous work on mixed-initiative dialogues focused on tracking a single thread of control among participants. We argue that this view of initiative fails to distinguish between *task initiative* and *dialogue initiative*, which together determine when and how an agent will address an issue. Although physical cues, such as gestures and eye gaze, play an important role in coordinating initiative shifts in face-to-face interactions, a great deal of information regarding initiative shifts can be extracted from utterances based on linguistic and domain knowledge alone. By taking into account such cues during dialogue interactions, the system is better able to determine

the task and dialogue initiative holders for each turn and to tailor its response to user utterances accordingly.

In this paper, we show how distinguishing between task and dialogue initiatives accounts for phenomena in collaborative dialogues that previous models were unable to explain. We show that a set of cues, which can be recognized based on linguistic and domain knowledge alone, can be utilized by a model for tracking initiative to predict the task and dialogue initiative holders with 99.1% and 87.8% accuracies, respectively, in collaborative planning dialogues. Furthermore, application of our model to dialogues in various other collaborative environments consistently increases the accuracies in the prediction of task and dialogue initiative holders by 2-4 and 8-13 percentage points, respectively, compared to a simple prediction method without the use of cues, thus illustrating the generality of our model.

## 2 Task Initiative vs. Dialogue Initiative

### 2.1 Motivation

Previous work on mixed-initiative dialogues focused on tracking and allocating a single thread of control, the *conversational lead*, among participants. Novick (1988) developed a computational model that utilizes *meta-locutionary acts*, such as *repeat* and *give-turn*, to capture mixed-initiative behavior in dialogues. Whittaker and Stenton (1988) devised rules for allocating dialogue control based on utterance types, and Walker and Whittaker (1990) utilized these rules for an analytical study on discourse segmentation. Kitano and Van Ess-Dykema (1991) developed a plan-based dialogue understanding model that tracks the conversational initiative based on the domain and discourse plans behind the utterances. Smith and Hipp (1994) developed a dialogue system that varies its responses to user utterances based on four dialogue modes which model different levels of initiative exhibited by dialogue participants. However, the dialogue mode is determined at the outset and cannot be changed during the dialogue. Guinn (1996) subsequently developed a system that allows change in the level of ini-

tiative based on initiative-changing utterances and each agent’s competency in completing the current subtask.

However, we contend that merely maintaining the conversational lead is insufficient for modeling complex behavior commonly found in naturally-occurring collaborative dialogues (SRI Transcripts, 1992; Gross, Allen, and Traum, 1993; Heeman and Allen, 1995). For instance, consider the alternative responses in utterances (3a)-(3c), given by an advisor to a student’s question:

- (1) *S: I want to take NLP to satisfy my seminar course requirement.*
- (2) *Who is teaching NLP?*
- (3a) *A: Dr. Smith is teaching NLP.*
- (3b) *A: You can’t take NLP because you haven’t taken AI, which is a prerequisite for NLP.*
- (3c) *A: You can’t take NLP because you haven’t taken AI, which is a prerequisite for NLP. You should take distributed programming to satisfy your requirement, and sign up as a listener for NLP.*

Suppose we adopt a model that maintains a single thread of control, such as that of (Whittaker and Stenton, 1988). In utterance (3a), A directly responds to S’s question; thus the conversational lead remains with S. On the other hand, in (3b) and (3c), A takes the lead by initiating a subdialogue to correct S’s invalid proposal. However, existing models cannot explain the difference in the two responses, namely that in (3c), A actively participates in the planning process by explicitly proposing domain actions, whereas in (3b), she merely conveys the invalidity of S’s proposal. Based on this observation, we argue that it is necessary to distinguish between *task initiative*, which tracks the lead in the development of the agents’ plan, and *dialogue initiative*, which tracks the lead in determining the current discourse focus (Chu-Carroll and Brown, 1997).<sup>1</sup> This distinction then allows us to explain A’s behavior from a response generation point of view: in (3b), A responds to S’s proposal by merely taking over the dialogue initiative, i.e., informing S of the invalidity of the proposal, while in (3c), A responds by taking over both the task and dialogue initiatives, i.e., informing S of the invalidity and suggesting a possible remedy.

An agent is said to have the *task initiative* if she is directing how the agents’ task should be accomplished, i.e., if her utterances directly propose *actions* that the

<sup>1</sup>Although independently conceived, this distinction between task and dialogue initiatives is similar to the notion of *choice of task* and *choice of speaker* in initiative in (Novick and Sutton, 1997), and the distinction between *control* and *initiative* in (Jordan and Di Eugenio, 1997).

	TI: system	TI: manager
DI: system	37 (3.5%)	274 (26.3%)
DI: manager	4 (0.4%)	727 (69.8%)

Table 1: Distribution of Task and Dialogue Initiatives

agents should perform. The utterances may propose *domain* actions (Litman and Allen, 1987) that directly contribute to achieving the agents’ goal, such as “*Let’s send engine E2 to Corning.*” On the other hand, they may propose *problem-solving* actions (Allen, 1991; Lambert and Carberry, 1991; Ramshaw, 1991) that contribute not directly to the agents’ domain goal, but to how they would go about achieving this goal, such as “*Let’s look at the first [problem] first.*” An agent is said to have the *dialogue initiative* if she takes the conversational lead in order to establish mutual beliefs, such as mutual beliefs about a piece of domain knowledge or about the validity of a proposal, between the agents. For instance, in responding to agent A’s proposal of sending a boxcar to Corning via Dansville, agent B may take over the dialogue initiative (but not the task initiative) by saying “*We can’t go by Dansville because we’ve got Engine 1 going on that track.*” Thus, when an agent takes over the task initiative, she also takes over the dialogue initiative, since a proposal of actions can be viewed as an attempt to establish the mutual belief that a set of actions be adopted. On the other hand, an agent may take over the dialogue initiative but not the task initiative, as in (3b) above.

## 2.2 An Analysis of the TRAINS91 Dialogues

To analyze the distribution of task/dialogue initiatives in collaborative planning dialogues, we annotated the TRAINS91 dialogues (Gross, Allen, and Traum, 1993) as follows: each dialogue turn is given two labels, *task initiative* (TI) and *dialogue initiative* (DI), each of which can be assigned one of two values, *system* or *manager*, depending on which agent holds the task/dialogue initiative during that turn.<sup>2</sup>

Table 1 shows the distribution of task and dialogue initiatives in the TRAINS91 dialogues. It shows that while in the majority of turns, the task and dialogue initiatives are held by the same agent, in approximately 1/4 of the turns, the agents’ behavior can be better accounted for by tracking the two types of initiatives separately.

To assess the reliability of our annotations, approximately 10% of the dialogues were annotated by two additional coders. We then used the kappa statistic (Siegel and Castellan, 1988; Carletta, 1996) to assess the level of agreement between the three coders with respect to the

<sup>2</sup>An agent holds the task initiative during a turn as long as *some* utterance during the turn directly proposes how the agents should accomplish their goal, as in utterance (3c).

task and dialogue initiative holders. In this experiment,  $K$  is 0.57 for the task initiative holder agreement and  $K$  is 0.69 for the dialogue initiative holder agreement.

Carletta suggests that content analysis researchers consider  $K > .8$  as good reliability, with  $.67 < K < .8$  allowing tentative conclusions to be drawn (Carletta, 1996). Strictly based on this metric, our results indicate that the three coders have a reasonable level of agreement with respect to the dialogue initiative holders, but do not have reliable agreement with respect to the task initiative holders. However, the kappa statistic is known to be highly problematic in measuring inter-coder reliability when the likelihood of one category being chosen overwhelms that of the other (Grove et al., 1981), which is the case for the task initiative distribution in the TRAINS91 corpus, as shown in Table 1. Furthermore, as will be shown in Table 4, Section 4, the task and dialogue initiative distributions in TRAINS91 are not at all representative of collaborative dialogues. We expect that by taking a sample of dialogues whose task/dialogue initiative distributions are more representative of all dialogues, we will lower the value of  $P(E)$ , the probability of chance agreement, and thus obtain a higher kappa coefficient of agreement. However, we leave selecting and annotating such a subset of representative dialogues for future work.

### 3 A Model for Tracking Initiative

Our analysis shows that the task and dialogue initiatives shift between the participants during the course of a dialogue. We contend that it is important for the agents to take into account signals for such initiative shifts for two reasons. First, recognizing and providing signals for initiative shifts allow the agents to better coordinate their actions, thus leading to more coherent and cooperative dialogues. Second, by determining whether or not it should hold the task and/or dialogue initiatives when responding to user utterances, a dialogue system is able to tailor its responses based on the distribution of initiatives, as illustrated by the previous dialogue (Chu-Carroll and Brown, 1997). This section describes our model for tracking initiative using cues identified from the user's utterances.

Our model maintains, for each agent, a *task initiative index* and a *dialogue initiative index* which measure the amount of evidence available to support the agent holding the task and dialogue initiatives, respectively. After each turn, new initiative indices are calculated based on the current indices and the effects of the cues observed during the turn. These cues may be explicit requests by the speaker to give up his initiative, or implicit cues such as ambiguous proposals. The new initiative indices then determine the initiative holders for the next turn.

We adopt the Dempster-Shafer theory of evidence (Shafer, 1976; Gordon and Shortliffe, 1984) as our un-

derlying model for inferring the accumulated effect of multiple cues on determining the initiative indices. The Dempster-Shafer theory is a mathematical theory for reasoning under uncertainty which operates over a set of possible outcomes,  $\Theta$ . Associated with each piece of evidence that may provide support for the possible outcomes is a *basic probability assignment (bpa)*, a function that represents the impact of the piece of evidence on the subsets of  $\Theta$ . A bpa assigns a number in the range  $[0,1]$  to each subset of  $\Theta$  such that the numbers sum to 1. The number assigned to the subset  $\Theta_1$  then denotes the amount of support the evidence directly provides for the conclusions represented by  $\Theta_1$ . When multiple pieces of evidence are present, Dempster's combination rule is used to compute a new bpa from the individual bpa's to represent their cumulative effect.

The reasons for selecting the Dempster-Shafer theory as the basis for our model are twofold. First, unlike the Bayesian model, it does not require a complete set of *a priori* and conditional probabilities, which is difficult to obtain for sparse pieces of evidence. Second, the Dempster-Shafer theory distinguishes between situations in which no evidence is available to support any conclusion and those in which equal evidence is available to support each conclusion. Thus the outcome of the model more accurately represents the *amount* of evidence available to support a particular conclusion, i.e., the *provability* of the conclusion (Pearl, 1990).

#### 3.1 Cues for Tracking Initiative

In order to utilize the Dempster-Shafer theory for modeling initiative, we must first identify the cues that provide evidence for initiative shifts. Whittaker, Stenton, and Walker (Whittaker and Stenton, 1988; Walker and Whittaker, 1990) have previously identified a set of utterance intentions that serve as cues to indicate shifts or lack of shifts in initiative, such as prompts and questions. We analyzed our annotated TRAINS91 corpus and identified additional cues that may have contributed to the shift or lack of shift in task/dialogue initiatives during the interactions. This results in eight cue types, which are grouped into three classes, based on the kind of knowledge needed to recognize them. Table 2 shows the three classes, the eight cue types, their subtypes if any, whether a cue may affect merely the dialogue initiative or both the task and dialogue initiatives, and the agent expected to hold the initiative in the next turn.

The first cue class, *explicit cues*, includes explicit requests by the speaker to give up or take over the initiative. For instance, the utterance “Any suggestions?” indicates the speaker's intention for the hearer to take over both the task and dialogue initiatives. Such explicit cues can be recognized by inferring the discourse and/or problem-solving intentions conveyed by the speaker's utterances.

Class	Cue Type	Subtype	Effect	Initiative	Example
Explicit	Explicit requests	give up	both	hearer	"Any suggestions?" "Summarize the plan up to this point"
		take over	both	speaker	"Let me handle this one."
Discourse	End silence		both	hearer	
	No new info	repetitions	both	hearer	A: "Grab the tanker, pick up oranges, go to Elmira, make them into orange juice." B: "We go to Elmira, we make orange juice, okay."
		prompts	both	hearer	"Yeah", "Ok", "Right"
	Questions	domain	DI	speaker	"How far is it from Bath to Corning?"
		evaluation	DI	hearer	"Can we do the route the banana guy isn't doing?"
	Obligation fulfilled	task	both	hearer	A: "Any suggestions?" B: "Well, there's a boxcar at Dansville." "But you have to change your banana plan." A: "How long is it from Dansville to Corning?"
		discourse	DI	hearer	A: "Go ahead and fill up E1 with bananas." B: "Well, we have to get a boxcar." A: "Right, okay. It's shorter to Bath from Avon."
Analytical	Invalidity	action	both	hearer	A: "Let's get the tanker car to Elmira and fill it with OJ." B: "You need to get oranges to the OJ factory."
		belief	DI	hearer	A: "It's shorter to Bath from Avon." B: "It's shorter to Dansville." "The map is slightly misleading."
	Suboptimality		both	hearer	A: "Using Saudi on Thursday the eleventh." B: "It's sold out." A: "Is Friday open?" B: "Economy on Pan Am is open on Thursday."
	Ambiguity	action	both	hearer	A: "Take one of the engines from Corning." B: "Let's say engine E2."
		belief	DI	hearer	A: "We would get back to Corning at 4." B: "4PM? 4AM?"

Table 2: Cues for Modeling Initiative

The second cue class, *discourse cues*, includes cues that can be recognized using linguistic and discourse information, such as from the surface form of an utterance, or from the discourse relationship between the current and prior utterances. It consists of four cue types. The first type is perceptible silence at the end of an utterance, which suggests that the speaker has nothing more to say and may intend to give up her initiative. The second type includes utterances that do not contribute information that has not been conveyed earlier in the dialogue. It can be further classified into two groups: *repetitions*, a subset of the *informationally redundant utterances* (Walker, 1992), in which the speaker paraphrases an utterance by the hearer or repeats the utterance verbatim, and *prompts*, in which the speaker merely acknowledges the hearer's previous utterance(s). Repetitions and prompts also suggest that the speaker has nothing more to say and indicate that the hearer should take over the initiative (Whittaker and Stenton, 1988). The third type includes questions which, based on anticipated responses, are divided into *domain* and *evaluation* questions. *Domain* questions are questions in which the speaker intends to obtain or verify a piece of domain knowledge. They usually merely require a direct response and thus typically do not result in an initiative shift. *Evaluation*

questions, on the other hand, are questions in which the speaker intends to assess the quality of a proposed plan. They often require an analysis of the proposal, and thus frequently result in a shift in dialogue initiative. The final type includes utterances that satisfy an outstanding task or discourse obligation. Such obligations may have resulted from a prior request by the hearer, or from an interruption initiated by the speaker himself. In either case, when the task/dialogue obligation is fulfilled, the initiative may be reverted back to the hearer who held the initiative prior to the request or interruption.

The third cue class, *analytical cues*, includes cues that cannot be recognized without the hearer performing an evaluation on the speaker's proposal using the hearer's private knowledge (Chu-Carroll and Carberry, 1994; Chu-Carroll and Carberry, 1995). After the evaluation, the hearer may find the proposal *invalid*, *suboptimal*, or *ambiguous*. As a result, he may initiate a subdialogue to resolve the problem, resulting in a shift in task/dialogue initiatives.<sup>3</sup>

<sup>3</sup>Whittaker, Stenton, and Walker treat subdialogues initiated as a result of these cues as interruptions, motivated by their collaborative planning principles (Whittaker and Stenton, 1988; Walker and Whittaker, 1990).

### 3.2 Utilizing the Dempster-Shafer Theory

As discussed earlier, at the end of each turn, new task/dialogue initiative indices are computed based on the current indices and the effect of the observed cues to determine the next task/dialogue initiative holders. In terms of the Dempster-Shafer theory, new task/dialogue bpa's ( $m_{t-new}/m_{d-new}$ )<sup>4</sup> are computed by applying Dempster's combination rule to the bpa's representing the current initiative indices<sup>5</sup> and the bpa of each observed cue.

Evidently, some cues provide stronger evidence for an initiative shift than others. Furthermore, a cue may provide stronger support for a shift in dialogue initiative than in task initiative. Thus, we associate with each cue two bpa's to represent its effect on changing the current task and dialogue initiative indices, respectively. We extended our annotations of the TRAINS91 dialogues to include, in addition to the agent(s) holding the task and dialogue initiatives for each turn, a list of cues observed during that turn. Initially, each cue<sub>i</sub> is assigned the following bpa's:  $m_{t-i}(\Theta) = 1$  and  $m_{d-i}(\Theta) = 1$ , where  $\Theta = \{\text{speaker, hearer}\}$ . In other words, we assume that the cue has no effect on changing the current initiative indices. We then developed a training algorithm (**Train-bpa**, Figure 1) and applied it on the annotated data to obtain the final bpa's.

For each turn, the task and dialogue bpa's for each observed cue are used, along with the current initiative indices, to determine the new initiative indices (step 2). The **combine** function utilizes Dempster's combination rule to combine pairs of bpa's until a final bpa is obtained to represent the cumulative effect of the given bpa's. The resulting bpa's are then used to predict the task/dialogue initiative holders for the next turn (step 3). If this prediction disagrees with the actual value in the annotated data, **Adjust-bpa** is invoked to alter the bpa's for the observed cues, and **Reset-current-bpa** is invoked to adjust the current bpa's to reflect the actual initiative holder (step 4).

**Adjust-bpa** adjusts the bpa's for the observed cues in favor of the actual initiative holder. We developed three adjustment methods by varying the effect that a disagreement between the actual and predicted initiative holders will have on changing the bpa's for the observed cues. The first is *constant-increment* where each time a disagreement occurs, the value for the actual initiative holder in the bpa is incremented by a constant ( $\Delta$ ), while

<sup>4</sup>Bpa's are represented by functions whose names take the form of  $m_{sub}$ . The subscript *sub* may be *t-X* or *d-X*, indicating that the function represents the task or dialogue bpa under scenario X.

<sup>5</sup>The initiative indices are represented as bpa's. For instance, the current task initiative indices take the following form:  $m_{t-cur}(\text{speaker}) = x$  and  $m_{t-cur}(\text{hearer}) = 1 - x$ .

**Train-bpa**(annotated-data):

1.  $m_{t-cur} \leftarrow$  default task initiative indices  
 $m_{d-cur} \leftarrow$  default dialogue initiative indices  
 $cur\text{-}data \leftarrow \text{read}(\text{annotated-data})$   
 $cue\text{-}set \leftarrow$  cues in  $cur\text{-}data$
2. */\* compute new initiative indices \*/*  
 $m_{t-obs} \leftarrow$  task initiative bpa's for cues in  $cue\text{-}set$   
 $m_{d-obs} \leftarrow$  dialogue initiative bpa's for cues in  $cue\text{-}set$   
 $m_{t-new} \leftarrow \text{combine}(m_{t-cur}, m_{t-obs})$   
 $m_{d-new} \leftarrow \text{combine}(m_{d-cur}, m_{d-obs})$
3. */\* determine predicted next initiative holders \*/*  
 If  $m_{t-new}(\text{speaker}) \geq m_{t-new}(\text{hearer})$ ,  
      $t\text{-predicted} \leftarrow \text{speaker}$   
 Else,  $t\text{-predicted} \leftarrow \text{hearer}$   
 If  $m_{d-new}(\text{speaker}) \geq m_{d-new}(\text{hearer})$ ,  
      $d\text{-predicted} \leftarrow \text{speaker}$   
 Else,  $d\text{-predicted} \leftarrow \text{hearer}$
4. */\* find actual initiative holders and compare \*/*  
 $new\text{-}data \leftarrow \text{read}(\text{annotated-data})$   
 $t\text{-actual} \leftarrow$  actual task initiative holder in  $new\text{-}data$   
 $d\text{-actual} \leftarrow$  actual dialogue initiative holder in  $new\text{-}data$   
 If  $t\text{-predicted} \neq t\text{-actual}$ ,  
     **Adjust-bpa**( $cue\text{-}set, task$ )  
     **Reset-current-bpa**( $m_{t-cur}$ )  
 If  $d\text{-predicted} \neq d\text{-actual}$ ,  
     **Adjust-bpa**( $cue\text{-}set, dialogue$ )  
     **Reset-current-bpa**( $m_{d-cur}$ )
5. If end-of-dialogue, return  
 Else, */\* swap roles of speaker and hearer \*/*  
 $m_{t-cur}(\text{speaker}) \leftarrow m_{t-new}(\text{hearer})$   
 $m_{d-cur}(\text{speaker}) \leftarrow m_{d-new}(\text{hearer})$   
 $m_{t-cur}(\text{hearer}) \leftarrow m_{t-new}(\text{speaker})$   
 $m_{d-cur}(\text{hearer}) \leftarrow m_{d-new}(\text{speaker})$   
 $cue\text{-}set \leftarrow$  cues in  $new\text{-}data$   
 Goto step 2.

Figure 1: Training Algorithm for Determining BPA's

that for  $\Theta$  is decremented by  $\Delta$ . The second method, *constant-increment-with-counter*, associates with each bpa for each cue a counter which is incremented when a correct prediction is made, and decremented when an incorrect prediction is made. If the counter is negative, the *constant-increment* method is invoked, and the counter is reset to 0. This method ensures that a bpa will only be adjusted if it has no "credit" for correct predictions in the past. The third method, *variable-increment-with-counter*, is a variation of *constant-increment-with-counter*. However, instead of determining whether an adjustment is needed, the counter determines the amount to be adjusted. Each time the system makes an incorrect prediction, the value for the actual initiative holder is incremented by  $\Delta/2^{count+1}$ , and that for  $\Theta$  decremented

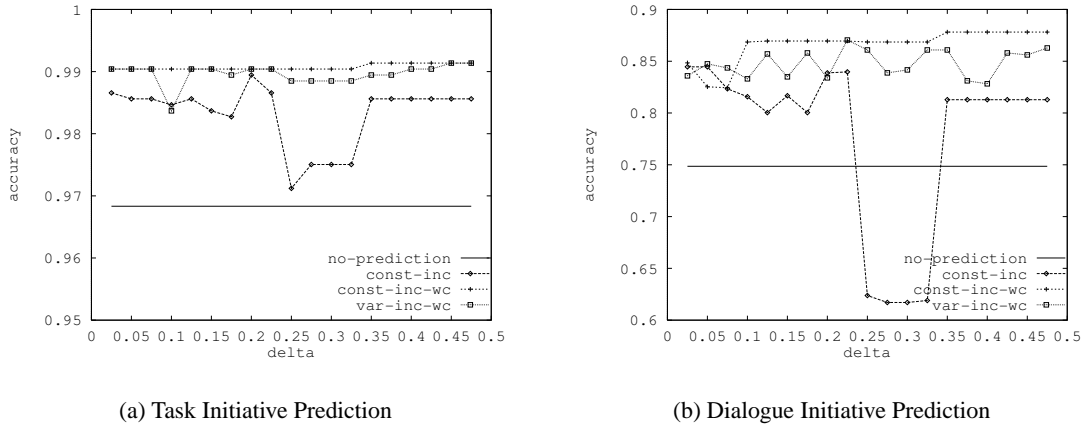


Figure 2: Comparison of Three Adjustment Methods

by the same amount.

In addition to experimenting with different adjustment methods, we also varied the increment constant,  $\Delta$ . For each adjustment method, we ran 19 training sessions with  $\Delta$  ranging from 0.025 to 0.475, incrementing by 0.025 between each session, and evaluated the system based on its accuracy in predicting the initiative holders for each turn. We divided the TRAINS91 corpus into eight sets based on speaker/hearer pairs. For each  $\Delta$ , we cross-validated the results by applying the training algorithm to seven dialogue sets and testing the resulting bpa's on the remaining set. Figures 2(a) and 2(b) show our system's performance in predicting the task and dialogue initiative holders, respectively, using the three adjustment methods.<sup>6</sup>

### 3.3 Discussion

Figure 2 shows that in the vast majority of cases, our prediction methods yield better results than making predictions without cues. Furthermore, substantial improvement is gained by the use of counters since they prevent the effect of the “exceptions of the rules” from accumulating and resulting in erroneous predictions. By restricting the increment to be inversely exponentially related to the “credit” the bpa had in making correct predictions, *variable-increment-with-counter* obtains better and more consistent results than *constant-increment*. However, the exceptions of the rules still resulted in undesirable effects, thus the further improved performance by *constant-increment-with-counter*.

We analyzed the cases in which the system, using

*constant-increment-with-counter* with  $\Delta = .35$ ,<sup>7</sup> made erroneous predictions. Tables 3(a) and 3(b) summarize the results of our analysis with respect to task and dialogue initiatives, respectively. For each cue type, we grouped the errors based on whether or not a shift occurred in the actual dialogue. For instance, the first row in Table 3(a) shows that when the cue *invalid action* is detected, the system failed to predict a task initiative shift in 2 out of 3 cases. On the other hand, it correctly predicted all 11 cases where no shift in task initiative occurred. Table 3(a) also shows that when an analytical cue is detected, the system correctly predicted all but one case in which there was no shift in task initiative. However, 55% of the time, the system failed to predict a shift in task initiative.<sup>8</sup> This suggests that other features need to be taken into account when evaluating user proposals in order to more accurately model initiative shifts resulting from such cues. Similar observations can be made about the errors in predicting dialogue initiative shifts when analytical cues are observed (Table 3(b)).

Table 3(b) shows that when a perceptible silence is detected at the end of an utterance, when the speaker utters a prompt, or when an outstanding discourse obligation is fulfilled (first three rows in table), the system correctly predicted the dialogue initiative holder in the vast majority of cases. However, for the cue class *questions*, when the actual initiative shift differs from the norm, i.e., speaker retaining initiative for evaluation questions and hearer taking over initiative for domain questions, the system's performance worsens. In the

<sup>7</sup>This is the value that yields the optimal results (Figure 2).

<sup>6</sup>For comparison purposes, the straight lines show the system's performance without the use of cues, i.e., always predict that the initiative remains with the current holder.

<sup>8</sup>In the case of suboptimal actions, we encounter the sparse data problem. Since there is only one instance of the cue in the set of dialogues, when the cue is present in the testing set, it is absent from the training set.

Cue Type	Subtype	Shift		No-Shift	
		error	total	error	total
Invalidity	action	2	3	0	11
Suboptimality		1	1	0	0
Ambiguity	action	3	7	1	5

(a) Task Initiative Errors

Cue Type	Subtype	Shift		No-Shift	
		error	total	error	total
End silence		13	41	0	53
No new info	prompts	7	193	1	6
Questions	domain	13	31	0	98
	evaluation	8	28	5	7
Obligation fulfilled	discourse	12	198	1	5
Invalidity		11	34	0	0
Suboptimality		1	1	0	0
Ambiguity		9	24	0	0

(b) Dialogue Initiative Errors

Table 3: Summary of Prediction Errors

case of domain questions, errors occur when 1) the response requires more reasoning than do typical domain questions, causing the hearer to take over the dialogue initiative, or 2) the hearer, instead of merely responding to the question, offers additional helpful information. In the case of evaluation questions, errors occur when 1) the result of the evaluation is readily available to the hearer, thus eliminating the need for an initiative shift, or 2) the hearer provides extra information. We believe that although it is difficult to predict when an agent may include extra information in response to a question, taking into account the cognitive load that a question places on the hearer may allow us to more accurately predict dialogue initiative shifts.

## 4 Applications in Other Environments

To investigate the generality of our system, we applied our training algorithm, using the *constant-increment-with-counter* adjustment method with  $\Delta = 0.35$ , on the TRAINS91 corpus to obtain a set of bpa's. We then evaluated the system on subsets of dialogues from four other corpora: the TRAINS93 dialogues (Heeman and Allen, 1995), airline reservation dialogues (SRI Transcripts, 1992), instruction-giving dialogues (Map Task Dialogues, 1996), and non-task-oriented dialogues (Switchboard Credit Card Corpus, 1992). In addition, we applied our baseline strategy which makes predictions without the use of cues to each corpus.

Table 4 shows a comparison between the dialogues

from the five corpora and the results of this evaluation. Row 1 in the table shows the number of turns where the *expert*<sup>9</sup> holds the task/dialogue initiative, with percentages shown in parentheses. This analysis shows that the distribution of initiatives varies quite significantly across corpora, with the distribution biased toward one agent in the TRAINS and maptask corpora, and split fairly evenly in the airline and switchboard dialogues. Row 2 shows the results of applying our baseline prediction method to the various corpora. The numbers shown are correct predictions in each instance, with the corresponding percentages shown in parentheses. These results indicate the difficulty of the prediction problem in each corpus that the task/dialogue initiative distribution (row 1) fails to convey. For instance, although the dialogue initiative is distributed approximately 30/70% between the two agents in the TRAINS91 corpus and 40/60% in the airline dialogues, the prediction rates in row 2 shows that in both cases, the distribution is the result of shifts in dialogue initiative in approximately 25% of the dialogue turns. Row 3 in the table shows the prediction results when applying our training algorithm using the *constant-increment-with-counter* method. Finally, the last row shows the improvement in percentage

<sup>9</sup>The *expert* is assigned as follows: in the TRAINS domain, the system; in the airline domain, the travel agent; in the map-task domain, the instruction giver; and in the switchboard dialogues, the agent who holds the dialogue initiative the majority of the time.

Corpus (# turns)	TRAINS91 (1042)		TRAINS93 (256)		Airline (332)		Maptask (320)		Switchboard (282)	
	task	dialogue	task	dialogue	task	dialogue	task	dialogue	task	dialogue
Expert control	41 (3.9%)	311 (29.8%)	37 (14.4%)	101 (39.5%)	194 (58.4%)	193 (58.1%)	320 (100%)	277 (86.6%)	N/A	166 (59.9%)
No cue	1009 (96.8%)	780 (74.9%)	239 (93.3%)	189 (73.8%)	308 (92.8%)	247 (74.4%)	320 (100%)	270 (84.4%)	N/A	193 (68.4%)
<i>const-inc-w-count</i>	1033 (99.1%)	915 (87.8%)	250 (97.7%)	217 (84.8%)	316 (95.2%)	281 (84.6%)	320 (100%)	297 (92.8%)	N/A	216 (76.6%)
<i>Improvement</i>	2.3%	12.9%	4.4%	11.0%	2.4%	10.2%	0.0%	8.4%	N/A	8.2%

Table 4: Comparison Across Different Application Environments

points between our prediction method and the baseline prediction method. To test the statistical significance of the differences between the results obtained by the two prediction algorithms, for each corpus, we applied Cochran’s  $Q$  test (Cochran, 1950) to the results in rows 2 and 3. The tests show that for all corpora, the differences between the two algorithms when predicting the task and dialogue initiative holders are statistically significant at the levels of  $p < 0.05$  and  $p < 10^{-5}$ , respectively.

Based on the results of our evaluation, we make the following observations. First, Table 4 illustrates the generality of our prediction mechanism. Although the system’s performance varies across environments, the use of cues consistently improves the system’s accuracies in predicting the task and dialogue initiative holders by 2-4 percentage points (with the exception of the maptask corpus in which there is no room for improvement)<sup>10</sup> and 8-13 percentage points, respectively. Second, Table 4 shows the specificity of the trained bpa’s with respect to application environments. Using our prediction mechanism, the system’s performances on the collaborative planning dialogues (TRAINS91, TRAINS93, and airline reservation) most closely resemble one another (last row in table). This suggests that the bpa’s may be somewhat sensitive to application environments since they may affect how agents interpret cues. Third, our prediction mechanism yields better results on task-oriented dialogues. This is because such dialogues are constrained by the goals; therefore, there are fewer digressions and offers of unsolicited opinion as compared to the switchboard corpus.

## 5 Conclusions

This paper discussed a model for tracking initiative between participants in mixed-initiative dialogue interactions. We showed that distinguishing between task and dialogue initiatives allows us to model phenomena in collaborative dialogues that existing systems are unable to explain. We presented eight types of cues that affect initiative shifts in dialogues, and showed how our model

predicts initiative shifts based on the current initiative holders and the effects that observed cues have on changing them. Our experiments show that by utilizing the *constant-increment-with-counter* adjustment method in determining the basic probability assignments for each cue, the system can correctly predict the task and dialogue initiative holders 99.1% and 87.8% of the time, respectively, in the TRAINS91 corpus, compared to 96.8% and 74.9% without the use of cues. The differences between these results are shown to be statistically significant using Cochran’s  $Q$  test. In addition, we demonstrated the generality of our model by applying it to dialogues in different application environments. The results indicate that although the basic probability assignments may be sensitive to application environments, the use of cues in the prediction process significantly improves the system’s performance.

## Acknowledgments

We would like to thank Lyn Walker, Diane Litman, Bob Carpenter, and Christer Samuelsson for their comments on earlier drafts of this paper, Bob Carpenter and Christer Samuelsson for participating in the coding reliability test, as well as Jan van Santen and Lyn Walker for discussions on statistical testing methods.

## References

- [Allen1991] Allen, James. 1991. Discourse structure in the TRAINS project. In *Darpa Speech and Natural Language Workshop*.
- [Carletta1996] Carletta, Jean. 1996. Assessing agreement on classification tasks: The kappa statistic. *Computational Linguistics*, 22:249–254.
- [Chu-Carroll and Brown1997] Chu-Carroll, Jennifer and Michael K. Brown. 1997. Initiative in collaborative interactions — its cues and effects. In *Working Notes of the AAAI-97 Spring Symposium on Computational Models for Mixed Initiative Interaction*, pages 16–22.
- [Chu-Carroll and Carberry1994] Chu-Carroll, Jennifer and Sandra Carberry. 1994. A plan-based model for

<sup>10</sup>In the maptask domain, the task initiative remains with one agent, the instruction giver, throughout the dialogue.



- response generation in collaborative task-oriented dialogues. In *Proceedings of the Twelfth National Conference on Artificial Intelligence*, pages 799–805.
- [Chu-Carroll and Carberry1995] Chu-Carroll, Jennifer and Sandra Carberry. 1995. Response generation in collaborative negotiation. In *Proceedings of the 33rd Annual Meeting of the Association for Computational Linguistics*, pages 136–143.
- [Cochran1950] Cochran, W. G. 1950. The comparison of percentages in matched samples. *Biometrika*, 37:256–266.
- [Gordon and Shortliffe1984] Gordon, Jean and Edward H. Shortliffe. 1984. The Dempster-Shafer theory of evidence. In Bruce Buchanan and Edward Shortliffe, editors, *Rule-Based Expert Systems: The MYCIN Experiments of the Stanford Heuristic Programming Project*. Addison-Wesley, chapter 13, pages 272–292.
- [Gross, Allen, and Traum1993] Gross, Derek, James F. Allen, and David R. Traum. 1993. The TRAINS 91 dialogues. Technical Report TN92-1, Department of Computer Science, University of Rochester.
- [Grove et al.1981] Grove, William M., Nancy C. Andreasen, Patricia McDonald-Scott, Martin B. Keller, and Robert W. Shapiro. 1981. Reliability studies of psychiatric diagnosis. *Archives of General Psychiatry*, 38:408–413.
- [Guinn1996] Guinn, Curry I. 1996. Mechanisms for mixed-initiative human-computer collaborative discourse. In *Proceedings of the 34th Annual Meeting of the Association for Computational Linguistics*, pages 278–285.
- [Heeman and Allen1995] Heeman, Peter A. and James F. Allen. 1995. The TRAINS 93 dialogues. Technical Report TN94-2, Department of Computer Science, University of Rochester.
- [Jordan and Di Eugenio1997] Jordan, Pamela W. and Barbara Di Eugenio. 1997. Control and initiative in collaborative problem solving dialogues. In *Working Notes of the AAAI-97 Spring Symposium on Computational Models for Mixed Initiative Interaction*, pages 81–84.
- [Kitano and Van Ess-Dykema1991] Kitano, Hiroaki and Carol Van Ess-Dykema. 1991. Toward a plan-based understanding model for mixed-initiative dialogues. In *Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics*, pages 25–32.
- [Lambert and Carberry1991] Lambert, Lynn and Sandra Carberry. 1991. A tripartite plan-based model of dialogue. In *Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics*, pages 47–54.
- [Litman and Allen1987] Litman, Diane and James Allen. 1987. A plan recognition model for subdialogues in conversation. *Cognitive Science*, 11:163–200.
- [Map Task Dialogues1996] Map Task Dialogues. 1996. Transcripts of DCIEM Sleep Deprivation Study, conducted by Defense and Civil Institute of Environmental Medicine, Canada, and Human Communication Research Centre, University of Edinburgh and University of Glasgow, UK. Distributed by HCRC and LDC.
- [Novick1988] Novick, David G. 1988. *Control of Mixed-Initiative Discourse Through Meta-Locutionary Acts: A Computational Model*. Ph.D. thesis, University of Oregon.
- [Novick and Sutton1997] Novick, David G. and Stephen Sutton. 1997. What is mixed-initiative interaction? In *Working Notes of the AAAI-97 Spring Symposium on Computational Models for Mixed Initiative Interaction*, pages 114–116.
- [Pearl1990] Pearl, Judea. 1990. Bayesian and belief-fusions formalisms for evidential reasoning: A conceptual analysis. In Glenn Shafer and Judea Pearl, editors, *Readings in Uncertain Reasoning*. Morgan Kaufmann, pages 540–574.
- [Ramshaw1991] Ramshaw, Lance A. 1991. A three-level model for plan exploration. In *Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics*, pages 36–46.
- [Shafer1976] Shafer, Glenn. 1976. *A Mathematical Theory of Evidence*. Princeton University Press.
- [Siegel and Castellan1988] Siegel, Sidney. and N. John. Castellan, Jr. 1988. *Nonparametric Statistics for the Behavioral Sciences*. McGraw Hill.
- [Smith and Hipp1994] Smith, Ronnie W. and D. Richard Hipp. 1994. *Spoken Natural Language Dialog Systems — A Practical Approach*. Oxford University Press.
- [SRI Transcripts1992] SRI Transcripts. 1992. Transcripts derived from audiotape conversations made at SRI International, Menlo Park, CA. Prepared by Jacqueline Kowtko under the direction of Patti Price.
- [Switchboard Credit Card Corpus1992] Switchboard Credit Card Corpus. 1992. Transcripts of telephone conversations on the topic of credit card use, collected at Texas Instruments. Produced by NIST, available through LDC.
- [Walker and Whittaker1990] Walker, Marilyn and Steve Whittaker. 1990. Mixed initiative in dialogue: An investigation into discourse segmentation. In *Proceedings of the 28th Annual Meeting of the Association for Computational Linguistics*, pages 70–78.

[Walker1992] Walker, Marilyn A. 1992. Redundancy in collaborative dialogue. In *Proceedings of the 15th International Conference on Computational Linguistics*, pages 345–351.

[Whittaker and Stenton1988] Whittaker, Steve and Phil Stenton. 1988. Cues and control in expert-client dialogues. In *Proceedings of the 26th Annual Meeting of the Association for Computational Linguistics*, pages 123–130.